

# Instruction Following in Text-Based Games

Mathieu Tuli, Andrew Li, Pashootan Vaezipoor, Toryn Q. Klassen<sup>†</sup>,  
Scott Sanner, Sheila A. McIlraith<sup>†</sup>

University of Toronto, Toronto, Canada

Vector Institute for Artificial Intelligence, Toronto, Canada

<sup>†</sup> Schwartz Reisman Institute for Technology and Society, Toronto, Canada

{mathieutuli, andrewli, pashootan, toryn, sheila}@cs.toronto.edu

ssanner@mie.utoronto.ca

## Abstract

Text-based games present a unique class of sequential decision making problem in which agents interact with a partially observable, simulated environment via actions and observations conveyed through natural language. Such observations typically include instructions that, in a reinforcement learning (RL) setting, can directly or indirectly guide a player towards completing reward-worthy tasks. In this work, we study the ability of RL agents to follow such instructions. We conduct experiments that show that the performance of state-of-the-art text-based game agents is largely unaffected by the presence or absence of such instructions, and that these agents are typically unable to execute tasks to completion. To further study and address the task of instruction following, we equip RL agents with an internal structured representation of natural language instructions in the form of Linear Temporal Logic (LTL), a formal language that is increasingly used for temporally extended reward specification in RL. Our framework both supports and highlights the benefit of understanding the temporal semantics of instructions and in measuring progress towards achievement of such a temporally extended behaviour. Experiments demonstrate the superior performance of our approach.

## 1 Introduction

Building AI agents that can understand natural language is an important and longstanding problem in AI. In recent years, instrumented text-based game (TBG) engines have served as compelling environments for studying a variety of tasks related to language understanding, affordance extraction, memory, and sequential decision making (e.g., (Côté et al., 2018; Adhikari et al., 2020; Liu et al., 2022)). They provide a simulated, partially observable environment where an agent can navigate and interact with environment objects, receiving

observations and administering commands via natural language. TextWorld (Côté et al., 2018) is a text-based game learning environment for training reinforcement learning (RL) agents. Successful play requires language understanding, effective navigation, memory, and an ability to follow instructions embedded within the text. Instructions may or may not be directly bound to reward but can guide an RL agent towards completing tasks and collecting reward.

In this paper we study instruction following in text-based games and propose an approach that advances the previous state of the art. To this end, we employ the state-of-the-art model-free TBG RL agent called GATA (Graph Aided Transformer Agent) (Adhikari et al., 2020) that operates in the TextWorld environment. GATA has made significant advances in performance by augmenting TBG agents with long-term memory – a critical component of effective game play. Despite GATA’s improvement over previous baselines, our experiments (see Figure 1) show that GATA performance is largely unaffected by the presence or absence of instructions, leading us to conclude that GATA is not effectively following instructions. We also find that while GATA agents are able to garner reward, they are not typically successful in *completing* tasks – an important vulnerability to the deployment of such techniques in environments where partial completion of tasks can be unsafe.

To further study and address the task of instruction following, we equip GATA with an internal structured representation of natural language instructions specified in Linear Temporal Logic (LTL) (Pnueli, 1977), a formal language that is increasingly used for temporally extended goals in planning and reward specification in RL. LTL also provides a mechanism to monitor progress towards completion of instructions. Our framework both supports and highlights the benefit of understanding the temporal semantics of instructions and in measuring progress towards achievement of a temporally extended behaviour. We perform experiments that illustrate the superior performance of our TBG agent and its ability to follow instructions. Contributions of this work include:

Figure 1: Comparison of GATA performance when trained with instructions (GATA<sup>D</sup>) versus when instructions are stripped from environment observations (GATA<sup>D-S</sup>). Agents were trained with 20 or 100 games, at increasing levels of task difficulty (level 1 vs level 2). Note that normalized game point performance (solid blocks) and rate of success (hashed blocks) are largely unchanged whether instructions are present or absent. Low success rate (i.e., task completion) rate is also seen in level 2.

- ^ Experiments that expose the lack of instruction following and low task completion rate in a state-of-the-art TBG agent.
- ^ An approach to the study and deployment of instruction following in TBG environments via exploitation of a formal language: LTL. LTL provides well-defined semantics and supports a measure of progress towards satisfaction of instructions.
- ^ An augmentation to an existing state-of-the-art architecture for TBGs to equip a TBG agent with instruction-following capabilities.
- ^ Comprehensive experiments and insights that study our and others' approaches to instruction following, and that highlight the superior performance of our proposed approach.

## 2 Background

In this section we introduce TextWorld, the TBG engine that we use, together with the Cooking domain that we employ in our experiments. We also overview Linear Temporal Logic, which (as described in [section 1](#)) we use in our approach as an internal representation for instructions.

### 2.1 Text-Based Games: TextWorld

Text-based games are partially observable multi-turn games where the environment and the player's action choices are represented textually. In this work, we use TextWorld ([Côté et al., 2018](#)) as our text-based game engine. A text-based game can be viewed as a (discrete-time) partially observable Markov decision process (POMDP)  $(\mathcal{S}; \mathcal{T}; \mathcal{A}; \mathcal{O}; \gamma; R; \rho)$  ([Côté et al., 2018](#)) where  $\mathcal{S}$  is the environment's state space,  $\mathcal{A}$  is the action space,  $T(s_{t+1} | s_t; a_t)$  where  $s_{t+1}, s_t \in \mathcal{S}$  and  $a_t \in \mathcal{A}$  is the conditional transition probability between states  $s_{t+1}$  and  $s_t$  given action  $a_t$ ,  $\mathcal{O}$  is the set of (partial) observations that the agent receives,  $(\alpha, j; s_t; a_{t-1})$

is the set of conditional observation probabilities,  $R: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward function, and  $\gamma \in [0, 1]$  is the discount factor. An agent's goal is to learn some optimal policy  $\pi^*(a|s)$  (or a policy that conditions on historical observations or on some internal memory) that maximizes the expected discounted return. In this work, we focus on the choice-based variant of games, similar to previous works ([Adhikari et al., 2020](#); [Narasimhan et al., 2015](#)). The action space  $\mathcal{A}$  is a list of possible commands and at each time-step  $t$  in the game, the agent must select action  $a_t \in \mathcal{C}_t$  from the current subset of permissible actions  $\mathcal{C}_t \subseteq \mathcal{A}$ .

#### 2.1.1 Environment Setting

We focus on the TextWorld Cooking domain popularized by ([Adhikari et al., 2020](#)) and Microsoft's First TextWorld Problems: A Language and Reinforcement Learning Challenge (FTWP) ([Trischler et al., 2019](#)). The game tasks agents with gathering and preparing various cooking ingredients described by an in-game recipe that is to be found. Game points (rewards) are earned for each of (1) collecting a required ingredient, (2) performing a preparatory step (some cutting or cooking action) on an ingredient as required by the recipe, (3) preparing the meal once all of the ingredients have been prepared, and (4) eating the meal. The game's partial observations can contain instructions that guide the agent towards completion of tasks, but not all instructions correspond directly to rewards. The game first instructs the agent to examine a cookbook, which elicits a recipe to be followed. The act of examining the cookbook returns no reward, but following its recipe will return reward. See [Appendix C](#) for more details. Success is determined by whether the recipe is fully completed and eaten. Preparing ingredients can also involve collecting certain tools (e.g., a knife). The game may also involve navigation—the agent may need to navigate to the kitchen or to find certain ingredients.

## 2.2 Linear Temporal Logic (LTL)

Linear Temporal Logic (LTL) (Pnueli, 1977) is a formal language a propositional logical language with temporal modalities that can be used to describe properties of trajectories. We will use LTL to specify instructions. LTL formulas are constructed from propositional variables (e.g., `player-has-carrot`), connectives from propositional logic (e.g., `!`), and two temporal operators: `(NEXT)` and `U (UNTIL)`. Formally, we define the syntax of LTL per (Baier and Katoen, 2008) as

$$\phi ::= p \mid \neg \phi \mid \phi \wedge \psi \mid \phi \vee \psi \mid \text{NEXT } \phi \mid \text{UNTIL } \phi \psi$$

where  $p \in P$  for some finite set of propositional symbols  $P$ . Satisfaction of an LTL formula is determined by a sequence of truth assignments  $\{h_0, h_1, h_2, \dots\}$  for  $P$ , where  $p \in P$  proposition  $p \in P$  holds at time step  $i$ . Formally,  $\{h_i\}$  satisfies  $\phi$  at time  $i = 0$ , denoted  $\{h_i\} \models \phi$ , under the following conditions:

$$\begin{aligned} \{h_i\} \models p & \iff p \in h_i, \text{ where } p \in P \\ \{h_i\} \models \neg \phi & \iff \{h_i\} \not\models \phi \\ \{h_i\} \models \phi \wedge \psi & \iff \{h_i\} \models \phi \text{ and } \{h_i\} \models \psi \\ \{h_i\} \models \phi \vee \psi & \iff \{h_i\} \models \phi \text{ or } \{h_i\} \models \psi \\ \{h_i\} \models \text{NEXT } \phi & \iff \{h_{i+1}\} \models \phi \\ \{h_i\} \models \text{UNTIL } \phi \psi & \iff \text{there exists } j \text{ such that} \\ & \{h_i\} \models \phi \text{ and } \{h_j\} \models \psi, \text{ and } \{h_k\} \models \phi \\ & \text{for all } k \in [i, j) \\ \{h_i\} \models \phi & \iff \{h_{i+1}\} \models \phi \end{aligned}$$

A sequence  $\{h_i\}$  is then said to satisfy  $\phi$  if  $\{h_0\} \models \phi$ .

Any LTL formula can be defined in terms of  $p \in P$ ,  $\neg$  (negation),  $\wedge$  (and), `(NEXT)`, and `U (UNTIL)`. From these operators, we can also define the Boolean operators `_` (or) and `!` (implication), and the temporal operators `(ALWAYS)` and `(EVENTUAL)`, where  $\{h_i\} \models \phi$  if  $\phi$  always holds in  $\{h_i\}$ , and  $\{h_i\} \models \phi$  if  $\phi$  holds at some point in  $\{h_i\}$ .

### 2.2.1 LTL Progression

LTL formulas can also be progressed along a sequence of truth assignments (Bacchus and Kabanza, 2000; Toro Icarte et al., 2018b) In other words, as an agent acts in the environment, resulting truth assignments can be used to update the formula to reflect what has been satisfied. The updated formula would now reflect the parts of the original formula that are remaining to be satisfied or whether the formula has been violated/satisfied. The progression operator `prog( ; ' )` is defined as:

Definition 2.1. For LTL formula  $\phi$ , truth assignment  $\{h_i\}$  over  $P$ , and  $p \in P$ , `prog( ; ' )` is defined as

$$\begin{aligned} \text{prog}( ; p) &= \begin{cases} \text{true} & \text{if } p \in h_i \\ \text{false} & \text{otherwise} \end{cases} \\ \text{prog}( ; \phi \wedge \psi) &= \text{prog}( ; \phi) \wedge \text{prog}( ; \psi) \\ \text{prog}( ; \neg \phi) &= \neg \text{prog}( ; \phi) \\ \text{prog}( ; \phi \text{ UNTIL } \psi) &= \text{prog}( ; \psi) \_ (\text{prog}( ; \phi) \wedge \phi \text{ UNTIL } \psi) \\ \text{prog}( ; \text{NEXT } \phi) &= \text{prog}( ; \phi) \end{aligned}$$

In the context of TextWorld, the progression operator can be applied at every step in the episode to update the LTL instruction fed to the agent. To do so, it's necessary to have event detectors that can detect when propositions are true as the agent acts during an episode (e.g., to detect that `player-has-carrot` is true when the player has the carrot). We discuss how event detection occurs in section 4. To illustrate how progression works, the LTL instruction `(EVENTUALLY player-has-carrot) ^ (EVENTUALLY player-has-apple)` would be progressed to `(EVENTUALLY player-has-apple)` once the agent grabs the carrot during an episode. In other words, once the player executes the `grab carrot` action, which results in a state where `player-has-carrot` is true, the LTL instruction is progressed to reflect that the agent no longer needs to get the carrot but must still grab the apple at some point.

## 3 Following Instructions in GATA

In order to evaluate the effectiveness of state-of-the-art text-based game agents at following instructions, we conducted experiments on the Cooking domain using the state-of-the-art model-free RL agent for TextWorld, GATA (Adhikari et al., 2020). GATA uses a transformer variant of the popular LSTM-DQN (Narasimhan et al., 2015) combined with a dynamic belief graph that is updated during game-play. The aim is to use this belief graph as long-term memory to improve action selection by modelling the underlying game dynamics (Adhikari et al., 2020). Formally, given the POMDP, GATA attempts to learn some optimal policy  $(\pi, g)$  where  $g$  is the belief graph.

While GATA's belief graph can capture goal relations (e.g. `apple-needs-cut`), it turns out that agents trained to condition on observations and the GATA belief graph alone largely ignore in-game instructions. We tested a GATA agent on levels 1 and 2 in the Cooking domain, after training on either the 20-game or 100-game training set, and found that in none of those settings was the cookbook examined more than 15% of the time (3/20 testing games). In short, the GATA agent usually doesn't observe what the recipe is for the current game meaning it has no way of knowing what the actual goal of the game is (except eventually from the rewards it gets and when the episode ends).

We further investigate how GATA agents fail to follow instructions by training these agents using modified game observations that have their instructions stripped (specifically, instructions directing the agent to examine the cookbook, the recipe text within the cookbook, and instructions to grab a knife if attempting to cut an ingredient without first holding the knife were removed from obser-

vations). This has two effects: (1) the agent no longer receives text-based instructions about what the goal is or what it should do; and (2) GATA's belief state will no longer capture goal relations like 'needs'. The results of this experiment are in [Figure 1](#), and demonstrate how GATA's performance remains largely unchanged. This suggests that GATA is (here at least) (a) not exploiting text-based instructions that would lead it to success and (b) even not exploiting the goal-related relations in its own belief state

The results in [Figure 1](#) also show a drop in GATA's performance when moving from level 1 to level 2 in the Cooking domain, where the games' complexity is increased by just one added ingredient preparation step in the recipes (see [Table 1](#) for more details on the levels). GATA has difficulty in fully completing tasks on level 2 games, where its success rate is roughly half that of its achieved normalized game points (only the latter metric was used by ([Adhikari et al., 2020](#))).

Given these insights, we wish to further study and address instruction following in TBGs. In the next section, we propose using LTL and demonstrate how existing work can be easily augmented.

## 4 An Approach to Following Instructions

We now investigate a mechanism for both studying and advancing the ability of an RL agent to follow instructions. We do so by translating instructions to an internal structured representation of language in the form of LTL, a formal language that is increasingly being used for reward specification in RL agents ([Vaezipoor et al., 2021](#); [Leon et al., 2020](#); [Kuo et al., 2020](#); [Camacho et al., 2019](#); [Toro Icarte et al., 2018b](#)). We describe how to augment the GATA architecture with these LTL instructions and how to monitor progress towards their completion.

### 4.1 Generating and Representing LTL Instructions for TextWorld

We use three types of instructions for the Cooking domain. The first instruction identifies the need to examine the cookbook: This instruction is defined as  $\phi$ : NEXT cookbook-is-examined: This instruction simply states that the agent should examine the cookbook (i.e. `cookbook-is-examined = true`) in the next step of the game. The second instruction is the actual recipe that gets elicited from the cookbook. We format this instruction to be order-invariant and incomplete. Order-invariance allows the agent to complete the instructions in any order, but is still constrained by any ordering that the TextWorld engine may enforce. Incomplete simply refers to the fact that not every single action required to complete the recipe is encoded (i.e. grabbing a knife before slicing a carrot, opening

the fridge). The agent must still learn to do these things to accomplish its tasks, but is not directly instructed to. Assuming the recipe requires that predicates  $p_1; p_2; \dots; p_n$  be true, the cookbook instructions are modelled as  $\phi : (\text{EVENTUALLY } p_1) \wedge (\text{EVENTUALLY } p_2) \wedge \dots \wedge (\text{EVENTUALLY } p_n)$ :

For example, in the Cooking Domain, this instruction might be the conjunction

$$\phi : (\text{EVENTUALLY } \text{apple-in-player}) \wedge (\text{EVENTUALLY } \text{meal-in-player}) \wedge (\text{EVENTUALLY } \text{meal-is-consumed}) :$$

The final type of instruction identifies the need to navigate to the kitchen. This instruction is defined as  $\phi$ : EVENTUALLY player-at-kitchen: This instruction will come prior to the first two described above, but is only used in games with navigation (see [Table 1](#)).

We build a simple LTL translator that generates these instructions from the textual observations, similar to the goal generator used in ([Liu et al., 2022](#)). TextWorld's observations are easily parsed to extract the goal information already contained within them, which we then formalize and keep track of using LTL. We provide examples of these observations and more details in [Appendix D](#). Note that these observations are only used to generate the instruction itself, and subsequently LTL progression is used with the GATA belief state as our event detector to monitor completion of instruction steps and to update instructions that remain to be addressed.

### 4.2 LTL Augmented Rewards and Episode Termination

We can also reward our agent for completing instructions, which we model as reward  $R(s; a; \phi)$ . For some labelling function  $L : S \rightarrow \mathcal{A} \rightarrow 2^P$  that assigns truth values to the propositions in  $P$ ,

$$R(s; a; \phi) = \begin{cases} 2 & \text{if } \text{prog}(L(s; a); \phi) = \text{true} \\ 1 & \text{if } \text{prog}(L(s; a); \phi) = \text{false} \\ 0 & \text{otherwise} \end{cases}$$

In other words, a bonus reward is given for every LTL instruction the agent satisfies and a penalty is given if the agent fails to complete an instruction. We perform an ablation study on the effect of this reward in [subsection G.4.2](#). We henceforth refer to this modified reward function as the LTL reward. The maximum bonus reward an agent receives is either 2 if there is no navigation task, or 3.

Further, because we wish to satisfy instructions, we can also use the instructions to modify episode termination. That is, if our LTL instruction is violated, we have arrived in a terminal state, even

if TextWorld has not indicated so. We perform an ablative study on the effect of this LTL-based termination in [subsection G.4.2](#).

### 4.3 LTL-GATA Model Architecture

We build a similar model to GATA's original architecture, augmented to include the LTL encoding of instructions and their progression according to observed system state. We dub this model LTL-GATA, which we describe in detail below. [Figure 2](#) depicts an episode step interaction of LTL-GATA with TextWorld and [Figure 3](#) depicts the model itself. Additional details can be found in [Appendix E](#).

**Graph Updater:** We use the original GATA-GTP model ([Adhikari et al., 2020](#)), which generates a discrete belief graph as a list of triplets of the form (object; relationship; object). It is composed of two sub-components: (a) the belief state updater, which generates  $g_t$  from observation  $o_t$  and the graph  $g_{t-1}$ ; and (b) the graph encoder, which encodes the current graph into a vector as  $G\mathbb{E}(g_t) = g_t^0 \in \mathbb{R}^D$  for some latent dimension  $D$ . The graph encoder is a relational graph convolutional network (R-GCN) ([Schlichtkrull et al., 2018](#)) using basis regularization ([Schlichtkrull et al., 2018](#)) and highway connections ([Srivastava et al., 2015](#)). We refer the readers to ([Adhikari et al., 2020](#)) for more details.

**LTL Updater:** The LTL updater generates and progresses LTL instructions. LTL instructions defining the need to arrive at the kitchen and examine the cookbook are generated from the initial observation  $o_0$ . The subsequent instruction defining the recipe is generated from game observation  $o_t$ , as described in [subsection 4.1](#), when the action examine cookbook is executed at time  $t$ . For the truth assignments (i.e. the labelling function  $L$ ), we leverage GATA's highly accurate belief state from the graph updater. We use the Spot engine ([Duret-Lutz et al., 2016](#)) to perform the progression.

**Text Encoders:** For encoding the action choices  $C_t$ , observations  $o_t$ , as well as encoding the LTL instructions  $'_t$ , we use a simplified version of the Transformer architecture presented in ([Vaswani et al., 2017](#)). This is the same architecture used in ([Adhikari et al., 2020](#)). For LTL instructions, we encode them directly as a string. For example, the LTL formula  $' : (\text{EVENTUAL}(p_1) \wedge (\text{EVENTUAL}(p_2)))$  where  $p_1 = \text{pepper-in-player}$  and  $p_2 = \text{pepper-is-cut}$ , has the string representation

```
str('') : \eventually player_has_pepper and
          eventually pepper_is_cut \
```

We format each predicate as a single token, and we show in [subsection G.4.1](#) that our method is robust to predicate format. For some input string  $v \in \mathbb{R}^D$  of length  $\ell$ , the text encoder outputs a single

vector  $TE(v) = v^0 \in \mathbb{R}^D$  of dimension  $D$ , which is the same latent dimension as the graph encoder.

**Action Selector:** The action selector is a 2-layer multi-layer perceptron (MLP). The encoded state vectors  $TE(o_t) = o_t^0 \in \mathbb{R}^D$ ,  $TE('_t) = '^0 \in \mathbb{R}^D$ , and  $G\mathbb{E}(g_t) = g_t^0 \in \mathbb{R}^D$  are concatenated to form the agent's final state representation  $z_t = [o_t^0; '^0; g_t^0] \in \mathbb{R}^{3D}$ . In contrast to ([Adhikari et al., 2020](#)), we concatenate features rather than use the bi-directional attention-based aggregator. This simplified the model's complexity and worked just as well experimentally. This vector is then repeated  $n_c$  times and concatenated with the encoded action choices  $C_t^0 \in \mathbb{R}^{n_c \times D}$  where  $n_c$  is the number of action choices. This input matrix is fed to the MLP which returns the a vector of Q-values for each action  $q_c \in \mathbb{R}^{n_c}$ .

**Training.** Formally, for belief state  $g$  and LTL instruction  $'$ , LTL-GATA aims to learn an optimal policy  $\pi(a|g; ')$ . To learn this optimal policy, we implement Double DQN (DDQN) ([Van Hasselt et al., 2016](#)) with reward function and termination criteria as discussed in [subsection 4.2](#). We use a prioritized experience replay buffer ([Schaul et al., 2016](#)). Refer to [subsection F.2](#) for further details.

## 5 Experiments

Our experimental assessment was designed both to understand how well GATA was exploiting observational instructions, as discussed in [section 3](#), and to assess the instruction-following performance of our proposed approach relative to this state of the art (not only in terms of game points but also successful completion). We additionally strove to assess features of our approach (such as monitoring instruction progress) that contributed to its performance, as well as general challenges to text-based game playing that limited its performance (such as navigation).

### 5.1 Experimental Setup

**Games.** Consistent with ([Adhikari et al., 2020](#)), we use the same set of games in the Cooking domain which has distinct training, validation, and testing sets. For the training games, there are two sets: one set that contains 20 unique games and another that contains 100 unique games. Both the validation and testing sets are set at 20 unique games each. The chosen levels are shown in [Table 1](#).

**Hyper-parameters.** We replicate all but three hyper-parameters from ([Adhikari et al., 2020](#)): (1) we use a batch size of 200 instead of 64 when training on the 100 game set, (2) for level 3, we use Boltzmann Action selection, and (3) we use Adam ([Kingma and Ba, 2015](#)) with a learning rate of 0.0003 instead of RAdam ([Liu et al., 2020](#)) with

<sup>1</sup><https://aka.ms/twkg/r1.0.2.zip>, accessed from <https://github.com/xingdi-eric-yuan/GATA-public>.

Figure 2: An example of a single step in an episode of TextWorld. The game environment returns an observation  $o_t$  and action candidate set  $C_t$  in response to action  $a_{t-1}$ . In turn, the agent's graph updater (GATA) updates its belief graph  $g_t$  in response to both  $o_t$  and  $g_{t-1}$ . Next,  $g_t$  and  $o_t$  update the LTL instructions.  $'_t$  is generated from  $o_t$  after the cookbook is examined and thereafter  $'_{t-1}$  is progressed to  $'_t$  at each time step. The policy network selects action  $a_t$  from  $C_t$  conditioned on  $o_t$ ,  $'_t$ , and  $g_t$  and the cycle repeats.

Table 1: Cooking Levels

Level	Recipe Size	Rooms	Max Score	Need {Grab, Cut, Cook}
0	1	1	3	{3, 7, 7}
1	1	1	4	{3, 3, 7}
2	1	1	5	{3, 3, 3}
3	1	9	3	{3, 7, 7}

Figure 3: LTL-GATA's policy model. The model chooses action  $a_t \in C_t$  conditioned on the state  $z_t = [o_t^0; 't; g_t^0]$ . The action selector chooses  $a_t$  based on the predicted Q-values.

a learning rate of 0.001. These changes boosted performance for all models. See subsection G.1 for more details.

Baselines. We compare against (1) TDQN (Adhikari et al., 2020), the transformer variant of the LSTM-DQN (Narasimhan et al., 2015) model, (2) GATA<sup>C</sup>, and (3) GATA<sup>D</sup>. GATA<sup>C</sup> is GATA's best performing model (GATA-COC) that uses a continue graph-updater pre-trained using contrastive observation classification. GATA<sup>D</sup> is a similarly performant model (GATA-GTP) that uses a discrete graph-updater pre-trained with ground-truth graphs from the FTWP dataset. Finally, we note that we found a few issues with GATA's original code<sup>2</sup> and have since fixed them (see subsection G.5.1). For comparison, we include the original paper GATA models, labelled as GATA<sup>C</sup> and GATA<sup>D</sup>.

Runs and measuring performance. Each experiment is run using 3 seeds (123, 321, and 666). We select the top-performing models (per seed) on the validation set during training and apply those models on the test set and report the average performance scores, as well as error bars signifying the standard deviation. We measure performance using two metrics: normalized accumulated game points and game success rate. We report averaged results over 3 seeds for each experiment. Previous works only compared using the normalized accumulated

game points, however this may sometimes be misleading – an agent could get  $3/4 = 0.75$  points on all games but never actually succeed on any. In contrast, measuring the success rate alongside the normalized game points allows for a more complete analysis of the agent's ability to play and complete these games.

## 5.2 LTL-GATA Compared to Baselines

Consistently high performance with 20 training games. We see from Figure 4 that LTL-GATA exhibits consistently high performance across levels as compared to the baselines when trained on the 20 games set. In particular, LTL-GATA maintains its performance on level 2, where the game's slight increase in complexity causes large performance drops in other methods. Our agent can easily complete the added task and maintain similar performance to the previous level 1.

Large performance gains with 100 training games. We see from Figure 4 that LTL-GATA gains considerable performance when trained on 100 games. With the added games, our agent is exposed to more predicates and can now generalize better to the testing set. Future work may look at how to achieve this kind of generalization without having to expose our agent to more predicates.

Success rate and normalized game points. Looking at the performance of GATA on level 2, it becomes apparent why measuring the success is important. Although it achieves almost 0.4 normal-

<sup>2</sup><https://github.com/xingdi-eric-yuan/GATA-public>, released under the open-source MIT License.

Figure 4: Testing scores across various levels and on both the 20 (top) and 100 (bottom) game training sets. We select the top-performing models (per seed) on the validation set during training and apply those models on the test set and report the average scores, as well as error bars signifying the standard deviation.

(a) Forcing GATA to Examine the Cookbook

(b) Progression Ablation

Figure 5: (a) A comparison of GATA<sup>D</sup> performance when given the option to examine the cookbook vs. when it is required to examine the cookbook. (b) A comparison of LTL-GATA with (Prog) and without (NoProg) using LTL progression.

normalized points, the actual success rate is near 0 for original GATA models, and 60% of the normalized points for the fixed models average across both training sets. In contrast, LTL-GATA exhibits high normalized points and success rate, where the average success rate across both training sets is 82% of the normalized points.

Competitive performance on level 3. Level 3 introduces the added challenge of navigation. LTL-GATA outperforms GATA in this level as well, but not to the degree of previous levels. Inspecting testing trajectories it becomes evident that both LTL-GATA and GATA methods struggle with navigation in this level, and have difficulties even navigating to the kitchen in the first place.

### 5.3 Does LTL Progression Matter?

We show in Figure 5b that the use of progression is critical to performance, where LTL-GATA without progression incurs a large performance drop-off, dropping below the performance of the baselines as well. Without progression, the LTL instruction

will not reflect the changes incurred by the agent's actions. This appears to confuse the agent considerably, demonstrated by its performance drop-off.

### 5.4 Forcing GATA to Examine the Cookbook

Because LTL-GATA is always tasked with examining the cookbook, we question whether a similar tasking for GATA improves performance. We experiment with GATA<sup>D</sup> by forcing the agent to examine the cookbook on the first step of the episode. Forcing GATA to examine the cookbook will elicit goal relations like (apple,needs,cut) in the belief state. We show however in Figure 5a that GATA does not improve when being given the cookbook. This shows that GATA cannot make use of the information elicited from the cookbook, continuing to ignore important instructions. Even with the presence of goal relations in its belief state, GATA fails to properly attend to this information. This highlights the benefits of a formalized representation of instructions used by LTL-GATA.

## 6 Related Work

Text-based games. In this work we equip a text-based deep RL agent with formalized LTL instructions, building on previous works that employed belief graphs for solving text-based games. (Adhikari et al., 2020) focused on supervised (i.e. translation) and self-supervised learned mechanisms to construct such belief graphs, whereas (Ammanabrolu and Hausknecht, 2020; Yin and May, 2019b; Ammanabrolu and Riedl, 2019) employed rule-based methods. At a larger scope, there is a host of other works on playing text-based games using deep reinforcement learning (Hausknecht et al., 2020; Zahavy et al., 2018; Jain et al., 2020; Yin and May, 2019a). (Yuan et al., 2018) used count-based memory to shape the reward to improve in exploration and generalization in a simple domain. (Narasimhan et al., 2015; He et al., 2016) proposed variations of an LSTM-based model, which the TDQN model used in this work is built from. In just published work, (Liu et al., 2022) took a model-based approach, focusing on object-oriented dynamics. However, these works do not address the role and representation of instructions that defines our work. (Kimura et al., 2021) does employ a neuro-symbolic RL method using Logical Neural Networks. However, it does not focus on instructions, operates over all logical facts of the environment, and is applied to a simpler domain.

Following LTL instructions. (Vaezipoor et al., 2021) trained an RL agent to follow various LTL instructions in both discrete and continuous action-space visual environments. They used R-GCNs to learn representations of the LTL instructions and employed LTL progression. Their model showed good generalization performance on similar and much larger unseen instructions than those observed during training. However, they relied on ground-truth event detectors, while we use a learned event detector (GATA) and opt for training the LTL semantics end-to-end using a transformer rather than an R-GCN. Works using LTL for reward specification (Leon et al., 2020; Kuo et al., 2020; Camacho et al., 2019; Toro Icarte et al., 2018b; Littman et al., 2017) or advice (Toro Icarte et al., 2018a) in RL agents exist, however they do not focus on text-based environments, which are inherently partially observable. Our work further differs in that we model instructions dynamically on the environment's observations and use a learned event detector (GATA) rather than ground truth.

Following natural language instructions. Following instructions conveyed via natural language has been studied in a number of different guises within AI (e.g., Mei et al., 2016; Artzi and Zettlemoyer, 2013; Chen and Mooney, 2011; Anderson et al., 2018; Tellex et al., 2011; MacMahon et al., 2006). Compared to LTL, which offers struc-

tured syntax and formal semantics, the ambiguity of natural language can make learning for instruction following more challenging. The merits of LTL for conveying instructions has been recognized in the past, resulting in several efforts to convert natural language into LTL for this and related purposes (e.g., (Finucane et al., 2010; Chen, 2018; Patel et al., 2020, 2019; Hahn et al., 2022)).

## 7 Conclusion

We studied the ability of RL agents to follow instructions in text-based games using TextWorld. We conducted experiments to show how current state-of-the-art model-free agents largely fail to exploit instructions and do not typically complete prescribed tasks, and how LTL can be used to construct internal structured representations for state augmentation that result in large performance improvements and more reliable instruction following and task completion. Experiments showed that monitoring instruction progress was critical to these gains. Our method inherits limitations in dealing with navigation and unseen games from prior work, but these concerns are somewhat orthogonal to our focus on instruction following. These present important complementary directions for future work. Finally, we consider the broader impact of this work by relating to the critical need for good instruction following in safety-oriented domains such as autonomous transport or health care. Overall, we intend this paper to highlight the importance of studying instruction following in environments like TextWorld that act as a proxies to the general class of problems dealing with language understanding and human-machine interaction.

## 8 Acknowledgements

Thank you to the reviewers for the engaging and constructive questions and the suggestions. We gratefully acknowledge funding from the Natural Sciences and Engineering Research Council of Canada (NSERC), the Canada CIFAR AI Chairs Program, and Microsoft Research. Resources used in preparing this research were provided, in part, by the Province of Ontario, the Government of Canada through CIFAR, and companies sponsoring the Vector Institute for Artificial Intelligence ([www.vectorinstitute.ai/partners](http://www.vectorinstitute.ai/partners)). Finally, we thank the Schwartz Reisman Institute for Technology and Society for providing a rich multi-disciplinary research environment.

## References

Ashutosh Adhikari, Xingdi Yuan, Marc-Alexandre Côté, Mikuláš Zelinka, Marc-Antoine Rondeau, Romain Laroche, Pascal Poupart, Jian Tang,



- Adam Trischler, and Will Hamilton. 2020. Learning dynamic belief graphs to generalize on text-based games. *Advances in Neural Information Processing Systems (NeurIPS2020)*33:3045-3057.
- Prithviraj Ammanabrolu and Matthew J. Hausknecht. 2020. [Graph constrained reinforcement learning for natural language action spaces](#). In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net.
- Prithviraj Ammanabrolu, Liwei Jiang, Maarten Sap, Hannaneh Hajishirzi, and Yejin Choi. 2022. [Aligning to social norms and values in interactive narratives](#). CoRR, abs/2205.01975.
- Prithviraj Ammanabrolu and Mark O. Riedl. 2019. [Playing text-adventure games with graph-based deep reinforcement learning](#). In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers), pages 3557-3565. Association for Computational Linguistics.
- Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian D. Reid, Stephen Gould, and Anton van den Hengel. 2018. [Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments](#). In 2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018, pages 3674-3683. Computer Vision Foundation / IEEE Computer Society.
- Yoav Artzi and Luke Zettlemoyer. 2013. [Weakly supervised learning of semantic parsers for mapping instructions to actions](#). *Transactions of the Association for Computational Linguistics*, 1:49-62.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. [Layer normalization](#). CoRR, abs/1607.06450.
- Fahiem Bacchus and Froduald Kabanza. 2000. Using temporal logics to express search control knowledge for planning. *Artificial Intelligence*, 116(1-2):123-191.
- Christel Baier and Joost-Pieter Katoen. 2008. *Principles of Model Checking* MIT Press.
- Alberto Camacho, Rodrigo Toro Icarte, Toryn Q. Klassen, Richard Anthony Valenzano, and Sheila A. McIlraith. 2019. LTL and beyond: Formal languages for reward function specification in reinforcement learning. In Proceedings of the 28th International Joint Conference on Artificial Intelligence (IJCAI), pages 6065-6073.
- David Chen and Raymond Mooney. 2011. Learning to interpret natural language navigation instructions from observations. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 25, pages 859-865.
- Menghan Chen. 2018. [Translating natural language into linear temporal logic](#). Review of Undergraduate Computer Science (RUCS).
- Marc-Alexandre Côté, Ákos Kádár, Xingdi Yuan, Ben Kybartas, Tavian Barnes, Emery Fine, James Moore, Matthew J. Hausknecht, Layla El Asri, Mahmoud Adada, Wendy Tay, and Adam Trischler. 2018. [Textworld: A learning environment for text-based games](#). In Computer Games - 7th Workshop, CGW 2018, Held in Conjunction with the 27th International Conference on Artificial Intelligence, IJCAI 2018, Stockholm, Sweden, July 13, 2018, Revised Selected Papers volume 1017 of Communications in Computer and Information Science, pages 41-75. Springer.
- Alexandre Duret-Lutz, Alexandre Lewkowicz, Amaury Fauchille, Thibaud Michaud, Etienne Renault, and Laurent Xu. 2016. Spot 2.0 A Framework for LTL and !-Automata Manipulation. In Proceedings of the 14th International Symposium on Automated Technology for Verification and Analysis (ATVA), pages 122-129. Springer.
- Cameron Finucane, Gangyuan Jing, and Hadas Kress-Gazit. 2010. [LTLMoP: Experimenting with language, temporal logic and robot control](#). In 2010 IEEE/RSJ International Conference on Intelligent Robots and Systems, October 18-22, 2010, Taipei, Taiwan, pages 1988-1993. IEEE.
- Christopher Hahn, Frederik Schmitt, Julia J. Tillman, Niklas Metzger, Julian Siber, and Bernd Finkbeiner. 2022. [Formal specifications from natural language](#). CoRR, abs/2206.01962.
- Matthew Hausknecht, Prithviraj Ammanabrolu, Marc-Alexandre Côté, and Xingdi Yuan. 2020. Interactive action games: A colossal adventure. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 34, pages 7903-7910.
- Ji He, Jianshu Chen, Xiaodong He, Jianfeng Gao, Lihong Li, Li Deng, and Mari Ostendorf. 2016. [Deep reinforcement learning with a natural language action space](#). In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers The Association for Computer Linguistics.
- Vishal Jain, William Fedus, Hugo Larochelle, Doina Precup, and Marc G. Bellemare. 2020. Algorithmic improvements for deep reinforcement learning applied to interactive action. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 34, pages 4328-4336.

- Daiki Kimura, Masaki Ono, Subhajit Chaudhury, Ryosuke Kohita, Akifumi Wachi, Don Joven Agravante, Michiaki Tatsubori, Asim Munawar, and Alexander Gray. 2021. [Neuro-symbolic reinforcement learning with rst-order logic](#). In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021, pages 3505-3511. Association for Computational Linguistics.
- Diederik P. Kingma and Jimmy Ba. 2015. [Adam: A method for stochastic optimization](#). In 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings.
- Yen-Ling Kuo, Boris Katz, and Andrei Barbu. 2020. [Encoding formulas as deep networks: Reinforcement learning for zero-shot execution of LTL formulas](#). In IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2020, Las Vegas, NV, USA, October 24, 2020 - January 24, 2021, pages 5604-5610. IEEE.
- Borja G Leon, Murray Shanahan, and Francesco Belardinelli. 2020. [Systematic Generalisation through Task Temporal Logic and Deep Reinforcement Learning](#). arXiv preprint arXiv:2006.08767.
- Michael L. Littman, Ufuk Topcu, Jie Fu, Charles Lee Isbell Jr., Min Wen, and James MacGlashan. 2017. [Environment-independent task specifications via GLTL](#). CoRR, abs/1704.04341.
- Guiliang Liu, Ashutosh Adhikari, Amir-massoud Farahmand, and Pascal Poupart. 2022. [Learning object-oriented dynamics for planning from text](#). In ICLR 2022.
- Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. 2020. [On the variance of the adaptive learning rate and beyond](#). In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net.
- Matt MacMahon, Brian Stankiewicz, and Benjamin Kuipers. 2006. [Walk the talk: Connecting language, knowledge, and action in route instructions](#). In Proceedings, The Twenty-First National Conference on Artificial Intelligence and the Eighteenth Innovative Applications of Artificial Intelligence Conference, July 16-20, 2006, Boston, Massachusetts, USA, pages 1475-1482. AAAI Press.
- Hongyuan Mei, Mohit Bansal, and Matthew R. Walter. 2016. [Listen, attend, and walk: Neural mapping of navigational instructions to action sequences](#). In Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, February 12-17, 2016, Phoenix, Arizona, USA, pages 2772-2778. AAAI Press.
- Tomas Mikolov, Edouard Grave, Piotr Bojanowski, Christian Puhresch, and Armand Joulin. 2017. [Advances in pre-training distributed word representations](#). arXiv preprint arXiv:1712.09405.
- Karthik Narasimhan, Tejas D. Kulkarni, and Regina Barzilay. 2015. [Language understanding for text-based games using deep reinforcement learning](#). In Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015, pages 1-11. The Association for Computational Linguistics.
- Roma Patel, Ellie Pavlick, and Stefanie Tellex. 2020. [Grounding language to non-markovian tasks with no supervision of task specifications](#). In Robotics: Science and Systems XVI, Virtual Event / Corvallis, Oregon, USA, July 12-16, 2020.
- Roma Patel, Roma Pavlick, and Stefanie Tellex. 2019. [Learning to ground language to temporal logical form](#). In NAACL.
- Amir Pnueli. 1977. [The temporal logic of programs](#). In Proceedings of the 18th IEEE Symposium on Foundations of Computer Science (FOCS), pages 46-57. IEEE.
- Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. 2016. [Prioritized experience replay](#). In 4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings.
- Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne van den Berg, Ivan Titov, and Max Welling. 2018. [Modeling relational data with graph convolutional networks](#). In European semantic web conference, pages 593-607. Springer.
- Rupesh Kumar Srivastava, Klaus Greff, and Jürgen Schmidhuber. 2015. [Highway networks](#). arXiv preprint arXiv:1505.00387.
- Stefanie Tellex, Thomas Kollar, Steven Dickerson, Matthew R. Walter, Ashis Gopal Banerjee, Seth J. Teller, and Nicholas Roy. 2011. [Understanding natural language commands for robotic navigation and mobile manipulation](#). In Proceedings of the Twenty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2011, San Francisco, California, USA, August 7-11, 2011. AAAI Press.
- Rodrigo Toro Icarte, Toryn Q. Klassen, Richard Valenzano, and Sheila A. McIlraith. 2018a. [Advice-based Exploration in Model-based Reinforcement Learning](#). In Proceedings of the 31st Canadian Conference on Artificial Intelligence (CCAI), pages 72-83. Springer.
- Rodrigo Toro Icarte, Toryn Q. Klassen, Richard Valenzano, and Sheila A. McIlraith. 2018b. [Teaching multiple tasks to an RL agent using LTL](#). In Proceedings of the 17th International Conference on Autonomous Agents and Multi-Agent Systems, pages 452-461.

Adam Trischler, Marc-Alexandre Côté, and Pedro Lima. 2019. First TextWorld Problems, the competition: Using text-based games to advance capabilities of AI agents. Microsoft Research Blog. URL <https://www.microsoft.com/en-us/research/blog/rst-textworld-problems-the-competition-using-text-based-games-to-advance-capabilities-of-ai-agents/>.

Pashootan Vaezipoor, Andrew C Li, Rodrigo A Toro Icarte, and Sheila A Mcilraith. 2021. LTL2action: Generalizing LTL instructions for multi-task RL. In International Conference on Machine Learning, pages 10497–10508. PMLR.

Hado Van Hasselt, Arthur Guez, and David Silver. 2016. Deep reinforcement learning with Double Q-Learning. In Proceedings of the AAAI conference on Artificial Intelligence, volume 30, pages 2094–2100.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. Advances in neural information processing systems, 30.

Xusen Yin and Jonathan May. 2019a. Comprehensible context-driven text game playing. In 2019 IEEE Conference on Games (CoG) pages 1–8. IEEE.

Xusen Yin and Jonathan May. 2019b. Learn how to cook a new recipe in a new house: Using map familiarization, curriculum learning, and bandit feedback to learn families of text-based adventure games. arXiv preprint arXiv:1908.04777.

Xingdi Yuan, Marc-Alexandre Côté, Alessandro Sordani, Romain Laroche, Remi Tachet des Combes, Matthew J. Hausknecht, and Adam Trischler. 2018. [Counting to explore and generalize in text-based games](#). CoRR, abs/1806.11525.

Tom Zahavy, Matan Haroush, Nadav Merlis, Daniel J. Mankowitz, and Shie Mannor. 2018. [Learn what not to learn: Action elimination with deep reinforcement learning](#). In Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada, pages 3566–3577.

# Appendices

A	Reinforcement Learning	12
B	Partially Observed Reinforcement Learning	12
C	TextWorld: Cooking Domain	12
D	Generating LTL in TextWorld	12
E	Model	13
E.1	Text Encoder	13
E.2	Encoder Independence	13
E.3	Action Selector	14
F	Implementation Details	15
F.1	Augmenting GATA's Pre-Training Dataset	15
F.2	Training	15
G	Experiments	15
G.1	Hyper-Parameters	15
G.2	Computational Requirements	16
G.3	Placeholder Header to Match Paper References	16
G.4	Additional Results	16
G.5	Code	18
G.6	Training Curves	18
H	Broader Impact	18

## A Reinforcement Learning

Reinforcement Learning (RL) is the problem of training machine learning models to solve sequential decision making problems. By interacting with an environment, RL agents must learn optimal behaviours given the current state of their environment. If the environment is fully observable, we can frame it as a Markov Decision Process (MDP) modelled as  $(\mathcal{S}; \mathcal{A}; \mathcal{T}; \mathcal{R}; \gamma)$  where  $\mathcal{S}$  is the environment's state space,  $\mathcal{A}$  is the action space,  $\mathcal{T}(s_{t+1}|s_t; a_t)$  where  $s_{t+1}; s_t \in \mathcal{S}$  and  $a_t \in \mathcal{A}$  is the conditional transition probability between states  $s_{t+1}$  and  $s_t$  given action  $a_t$ ,  $r_t = R(s; a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{R}$  is the reward function for state action pair  $(s; a)$ , and  $\gamma \in [0; 1]$  is the discount factor. The goal for an RL agent is to learn some optimal policy  $(\pi)$  that maximizes the expected discounted return  $E \sum_{k=0}^{\infty} \gamma^k r_{t+k} | \mathcal{S}_t = s$ . A single game is an episode and steps in an episode are indexed by  $t$ .

## B Partially Observed Reinforcement Learning

In a partially observed environment, an agent does not have access to the full state space  $\mathcal{S}$ . We can

frame this environment as a Partially Observable MDP (POMDP) modelled by  $(\mathcal{h}\mathcal{S}; \mathcal{A}; \mathcal{T}; \mathcal{O}; \gamma; \mathcal{R}; \mathcal{i})$ . In this new setting,  $\mathcal{h}\mathcal{S}; \mathcal{A}; \mathcal{T}; \mathcal{R}; \mathcal{i}$  remain unchanged,  $\mathcal{O}$  represents the set of (partial) observations that the agent receives and  $(o_t | s_t; a_{t-1})$  is the set of conditional observation probabilities. An agent's goal is to learn some optimal policy  $(\pi)$  (or a policy that conditions on historical observations or on some internal memory) that maximizes the expected discounted return.

## C TextWorld: Cooking Domain

We present two examples of observations with instructions in Table 2 and highlight where the instructions are and where the rewards come from.

## D Generating LTL in TextWorld

We provide some examples of the LTL instructions used in this work in Table 3, Table 4, and Table 5. We build a simple translator that reads game observations and constructs these LTL instructions directly, but only once. Repeated observations will not result in the same LTL formula being generated. Once a formula has been generated, LTL progression is used with the agent's belief state to progress the instructions along the truth assignments: observations are not directly used in the progression, although they do indirectly affect the progression by affecting the belief state.

For levels 0, 1, and 2, the LTL instructions that an agent can receive throughout an episode are (a) the task to examine the cookbook and (b) the recipe-bound task. In other words, the set of unprogressed instructions it can receive over the course of an episode (assuming the cookbook is examined) is as follows:

```
[ NEXT cookbook-is-examined;
  (EVENTUALLY p1) ^ (EVENTUALLY p2) ^
  ... (EVENTUALLY pn)]
```

where the recipe requires that predicates  $p_1; p_2; \dots; p_n$  be true. Note that we also consider eating the meal to be a part of recipe in this case, although it is not explicitly mentioned in the recipe. Further, we note that the prepare meal task is represented by the predicate meal-in-player, as this is the event that occurs when the meal is prepared in the game.

For levels with navigation (i.e. level 3),

```
[ EVENTUALLY player-at-kitchen;
  NEXT cookbook-is-examined;
  (EVENTUALLY p1) ^ (EVENTUALLY p2) ^
  ... (EVENTUALLY pn)]
```

Table 2: TextWorld observations for the Cooking Domain game. We show the observations and highlight where the instructions are, and naturally identify what the rewards would be. This is for a level 2 game, and the total possible reward is 5.

Initial Game Observation	
<p>You are hungry! Let's cook a delicious meal. <b>Check the cookbook in the kitchen for the recipe.</b> <b>Once done, enjoy your meal!</b> == kitchen == you find yourself in a kitchen. You start to take note of what's in the room. You can make out a closed fridge nearby. You can see an oven. You can make out a table. You wonder idly who left that here. You see a knife on the table. Something scurries by right in the corner of your eye. Probably nothing. You see a counter. The counter is vast. On the counter you see a raw red potato and a cookbook. You see a stove, but the thing is empty, unfortunately.</p>	
Reward	
<p>There is a reward of 1 given for eating the meal. i.e. the instruction <b>Once done, enjoy your meal!</b> will result in a reward of 1 after the recipe has been completed. Note that the instruction <b>Check the cookbook in the kitchen for the recipe.</b> is not bound to a reward.</p>	
Observation following the examine cookbook action	
<p>You open the copy of <b>Cooking : a modern approach (3rd ed.)</b> and start reading: recipe #1 - <b>Gather all following ingredients and follow the directions to prepare this tasty meal.</b> <b>Ingredients:</b> red potato: <b>directions:</b> chop the red potato, fry the red potato, prepare meal</p>	
Reward	
<p>There are 4 rewards from the instruction <b>Gather all following ingredients and follow the directions to prepare this tasty meal.</b> <b>Ingredients:</b> red potato: <b>directions:</b> chop the red potato, fry the red potato, prepare meal :</p> <ul style="list-style-type: none"> <li>^ 1 for grabbing the red potato</li> <li>^ 1 for chopping the red potato</li> <li>^ 1 for frying the red potato</li> <li>^ 1 for preparing the meal</li> </ul>	

where the agent has the added task of first navigating to the kitchen. This instruction provides no help for actually how to arrive at the kitchen, only that the agent must do so. As a result, LTL-GATA still suffers from the difficulties of exploration, and perhaps investigating how LTL can be used to improve in navigation could be a direction for future work.

In total, LTL generation occurs only twice for any level, either during the initial observation or when the cookbook is read. When multiple instructions are generated at once, the agent will process them sequentially, in the order they are given.

## E Model

### E.1 Text Encoder

The text encoder is a simple transformer-based model, with a transformer block (Vaswani et al., 2017) and word embedding layer. We use the pre-trained 300-dimensional fastText (Mikolov et al., 2017) word embeddings, which are trained on Com-

mon Crawl (600B tokens). These word embeddings are frozen during training. Strings are tokenized by spaces.

The transformer block is composed of: (1) a stack of 5 convolutional layers, (2) a single-head self-attention layer, and (3) a 2-layer MLP with ReLU non-linear activation function in between. The convolutional layers each have 64 filters, with kernel sizes of 5 and are each followed by a Layer Norm (Ba et al., 2016). We also use standard positional encoding (Vaswani et al., 2017). The self-attention layer uses a hidden size of 64. The Text Encoder outputs a single feature vector  $v \in \mathbb{R}^D$ , where  $D = 64$  in our experiments.

### E.2 Encoder Independence

Figure 3 in the main paper visualizes each component of our model. Specifically, our model has four encoders: (1) Graph Encoder, (2) Text Encoder for observations, (3) Text Encoder for LTL instructions, and (4) Text Encoder for action choices. We note here that each of these encoders are in-

Table 3: Level 3 observation and resulting generated LTL instruction

Observation
You are hungry! Let's cook a delicious meal. Check the cookbook in the kitchen for the recipe. Once done, enjoy your meal! -= corridor -= You've entered a corridor. There is a closed screen door leading west. You don't like doors? Why not try going north, that entranceway is not blocked by one. You need an exit without a door? You should try going south.
Generated LTL
This observation will generate two instructions: First,
' : (EVENTUALLY player-at-kitchen )
and second,
' : (NEXT cookbook-is-examined)

Table 4: Level 1 observation and resulting generated LTL instruction

Observation
You open the copy of Cooking : a modern approach (3rd ed.) and start reading: recipe #1 - Gather all following ingredients and follow the directions to prepare this tasty meal. Ingredients: red potato: directions: chop the red potato, prepare meal
Generated LTL
' : (EVENTUALLY red-potato-in-player ) ^ (EVENTUALLY red-potato-is-chopped ) ^ (EVENTUALLY meal-in-player ) ^ (EVENTUALLY meal-is-consumed ) :

dependent models, trained concurrently. This is in contrast to the original GATA model that used the same Text Encoder for both the actions and the observations. Because these Text Encoders are relatively small transformers, there is no issues with fitting this model in memory. As shown in Table 6, the model is still quite efficient, even more than the original GATA code. We found that using independent encoders resulted in better performance than using a single Text Encoder that would have been responsible for encoding the observations, LTL instructions, and action choices.

### E.3 Action Selector

The action selector is a simple two-layer MLP with a ReLU non-linear activation function in between. It takes as input, at time step  $t$ , the concatenated representation of the agent's state vector  $z_t \in \mathbb{R}^{3D}$  and the action choices  $C_t^0 \in \mathbb{R}^{n_c \times D}$ . Recall that in our experiments  $D = 64$ . The first layer uses an input dimension of  $4D$  and an output dimension of  $D$ . The second layer has an input dimension of  $D$

and output dimension of 1, which after squeezing the last dimension during the forward pass, the final output vector  $q_t \in \mathbb{R}^{n_c}$  represents the  $q$ -values for each action choice.

The input to the action selector is constructed by repeating the agent's state representation,  $z_t$ ,  $n_c$  times and then concatenating with the encoded actions  $C_t^0$ . We wanted to further explain why this occurs, as it may not be immediately clear. The action selector in this work is a parameter-tied Q-value predictor. That is, for some action  $a_i \in C_t$ ;  $i \in [1; \dots; n_c]$  and agent state representation  $z_t$ , the predicted Q-value is  $q = AS[a_i; z_t]$ . Thus, the action selector (i.e.  $AS$ ) predicts Q-values given action  $a_i$  and agent state representation  $z_t$ . Thus, during a single episode step, given our encoded actions  $C_t^0 \in \mathbb{R}^{n_c \times D}$ , in order for the action selector to predict Q-values for each of these action choices, we repeat  $z_t \in \mathbb{R}^{3D}$   $n_c$  times and stack it together, which results in a state matrix  $Z_t \in \mathbb{R}^{n_c \times 3D}$ . When we concatenate this matrix with our action choices we are left with the input to our

Table 5: Level 2 observation and resulting generated LTL instruction

Observation
You open the copy of <i>Cooking : a modern approach</i> (3rd ed.) and start reading: recipe #1 - Gather all following ingredients and follow the directions to prepare this tasty meal. Ingredients: red potato: directions: chop the red potato, fry the red potato, prepare meal
Generated LTL
$\varphi : (\text{EVENTUALLY}_{\text{red-potato-in-player}}) \wedge (\text{EVENTUALLY}_{\text{red-potato-is-chopped}}) \wedge (\text{EVENTUALLY}_{\text{red-potato-is-fried}}) \wedge (\text{EVENTUALLY}_{\text{meal-in-player}}) \wedge (\text{EVENTUALLY}_{\text{meal-is-consumed}})$

action selector:  $[C_t^0; Z_t] \in \mathbb{R}^{n_c \times 4D}$ . Looking at this matrix, each row in this input matrix is effectively the concatenation of action  $a_i$  with agent state representation  $z_t$ , and so passing this matrix to our action selector performs the parameter-tied Q-value prediction  $q_i = A\mathcal{S}[a_i; z_t]$  for all action choices, and outputs a single vector of Q-values for each action  $q_t \in \mathbb{R}^{n_c}$ . We can then use these predicted Q-values to perform action selection using either a greedy approach, an  $\epsilon$ -greedy approach, Boltzmann action selection, etc.

## F Implementation Details

### F.1 Augmenting GATA's Pre-Training Dataset

We note here that although possible, the vocabulary and dataset used by (Adhikari et al., 2020) did not allow for the knowledge triplets `f cookbook; is; examined` to be extracted from observations. Without this triple being extracted and added to the agent's belief state, there would be no way for the agent to progress LTL instructions requiring the agent to examine the cookbook. In our pre-training of the GATA graph encoder, we augmented the dataset provided by (Adhikari et al., 2020) to include the triplet `f cookbook; is; examined` when relevant (i.e. when the agent examines the cookbook). This was a simple process of adding this triple to the ground truth belief graphs in the dataset so that during pre-training, GATA could learn how to translate these triplets from relevant observations.

### F.2 Training

For training to learn our optimal policy we use the Double-DQN (DDQN) (Van Hasselt et al., 2016) framework. We use  $\epsilon$ -greedy for training, which first starts with a warm-up period, using a completely random policy (i.e.  $\epsilon = 1.0$ ) for the first 1;000

episodes. We then anneal from 1.0 to 0.1 over the next 3;000 episodes after the initial warm-up (i.e. episodes 1;000 to 4;000). We use a prioritized experience replay buffer ( $\alpha = 0.6$  and  $\beta = 0.4$ ) with capacity 500,000. For DDQN, the target network updates occur every 500 episodes. We update network parameters every 50 game steps, and we play 50 games in parallel.

We train all agents for 100,000 episodes using a discount factor of 0.9, and we use 123;321;666 as our random seeds. Each episode during training is limited to a maximum of 50 steps, and during testing/validation this limit is increased to 100 steps. We report results and save checkpoints every 1;000 episodes. We also use a patience window that reloads from the previous best checkpoint during training when validation performance has decreased for  $p$  episodes in a row. This is the same strategy used in (Adhikari et al., 2020). For our experiments, we used  $p = 3$ .

For reporting testing results, each model is trained using the three seeds mentioned before, and re-tuned on the validation set. That is, the checkpoint of the model that performs best on the validation set during training is saved, and each of these models (three, one for each seed) is applied to the test set. Reported test results are the average over these three models.

## G Experiments

### G.1 Hyper-Parameters

To have as fair a comparison as possible, we replicate all but three hyper-parameters from the settings used in (Adhikari et al., 2020). We do this to remove any bias towards more finely tuned experimental configurations and focus only on the LTL integration. Further, we re-run the GATA experiments to confirm their original results. The three changes we implemented were (1) we use a

Table 6: Training times for each model and training set size. The times were reported using a workstation with dual RTX3090s, an AMD Ryzen 5950x 16-core CPU, and 128GB of RAM. For the graph updater, COC stands for the contrastive observation classification pre-training (the continuous belief graph model) and GTP stands for ground-truth pre-training (the discrete belief graph model).

Model	Training Set Size	Batch Size	Approximate Time
TDQN	20	64	16 hours
LTL-GATA	20	64	24 hours
GATA <sup>D</sup>	20	64	24 hours
GATA <sup>C</sup>	20	64	24 hours
GATA <sup>D<sub>P</sub></sup>	20	64	36 hours
GATA <sup>C<sub>P</sub></sup>	20	64	36 hours
TDQN	100	200	32 hours
LTL-GATA	100	200	48 hours
GATA <sup>D*</sup>	100	200	48 hours
GATA <sup>C</sup>	100	200	48 hours
GATA <sup>D<sub>P</sub></sup>	100	200	65 hours
GATA <sup>C<sub>P</sub></sup>	100	200	65 hours
Graph Updater using COC	N/A	64	48 hours
Graph Updater using GTP	N/A	64	48 hours

batch size of 200 instead of 64 when training on the 100 game set, (2) for level 3, we use Boltzmann Action selection, and (3) we use Adam (Kingma and Ba, 2015) with a learning rate of 0.0003 instead of RAdam (Liu et al., 2020) with a learning rate of 0.001. These changes boosted performance for all models. For the 20 training game set, we use a batch size of 64.

For Boltzmann action selection, we used a temperature of  $\tau = 100$ . We experimented with various temperatures ( $\tau \in \{1, 10, 25, 50, 100, 200\}$ ) and found  $\tau = 100$  to perform the best across models.

## G.2 Computational Requirements

We report the wall-clock times for our experiments in Table 6.

## G.3 Placeholder Header to Match Paper References

This section header is included to match references from the main paper, please ignore.

## G.4 Additional Results

### G.4.1 Ablation: Formatting LTL Predicates

As we saw from Figure 4, LTL-GATA when trained on the 100 games set performs significantly better than when trained on the 20 game set, which we attribute to the increased exposure to predicates during training, allowing it to generalize better during testing. To see if we can achieve the same level of generalization when training on the 20 game set, we compare LTL-GATA with LTL predicates represented as single tokens (what we did in the

main paper) with using multiple tokens. That is, we compare the following two string representations:

(single-token predicates)  $\text{str}(\cdot)$  :  
 eventually player\_has\_pepper and  
 eventually pepper\_is\_cut  
 (multi-token predicates)  $\text{str}(\cdot)$  :  
 eventually player has pepper and  
 eventually pepper is cut

The single-token predicates are mapped in the vocabulary to a single word embedding. In our work, we compute word embedding for these single-token predicates by averaging the word embeddings of each underscore-separated word in the predicate. For example, the word embedding  $WE$  of the token `player_has_pepper` is

$$WE(\text{player\_has\_pepper}) = \frac{WE(\text{player}) + WE(\text{has}) + WE(\text{pepper})}{3}$$

For multiple-token predicates, each word has its own word embedding and we treat each word as any other word in the sentence. The idea is that by separating the tokens in the predicates, the text encoder (transformer) may be able to attend to each token independently, and during testing have better generalization. We visualize the results of this study in Figure 6. We can see from Figure 6 that this in fact does not help, and LTL-GATA performs almost equally in either scenario. This does however show how our method is robust to predicate format.



Figure 6: Study on LTL predicate format with single-token (Single) predicates and multi-token (Multi) predicates. Performance is largely unchanged with predicate format. The numbers overlaid on the bars indicate the numerical value of the normalized points achieved.



Figure 7: Ablation studies on LTL-based episode termination and new reward function  $R(s; a; ')$  and on LTL progression. (a) LTL-GATA with new reward function  $R(s; a; ')$  and without LTL-based episode termination (LTL-GATA-X). (b) LTL-GATA with base game reward function  $R(s; a)$  and with LTL-based episode termination (LTL-GATA-B). (c) LTL-GATA with base game reward function  $R(s; a)$  and without LTL-based episode termination (LTL-GATA-BX). The numbers overlaid on the bars indicate the numerical value of the normalized points achieved.

#### G.4.2 The Effect of LTL Reward and LTL-Based Termination

It is important to study the effect that the additional LTL bonus reward and LTL-based episode termination has on the performance of LTL-GATA. To study this, we consider three scenarios (a) LTL-GATA with the new reward function  $R(s; a; ')$  and without LTL-based episode termination; (b) LTL-GATA with the base TextWorld reward function  $R(s; a)$  and LTL-based episode termination; and (c) LTL-GATA with the normal TextWorld reward function  $R(s; a)$  and without LTL-based episode termination. For (a) and (b) we select level 1 on both the 20 and 100 game training set and level 2 on the 20 game training set. For (c) we select level 2 on 20 training games. We visualize the ablative study of these three scenarios in Figure 7.

From Figure 7 we can conclude that the presence of either the new reward function  $R(s; a; ')$  or LTL-based episode termination is important to the performance of LTL-GATA. This is because either of these methods will incentivize the agent to complete the initial NEXT cookbook-is-examined instruction, which isn't intrinsically rewarded by TextWorld. We can demonstrate the importance of this incentive by analyzing just one level (level 2 on 20 training games). Removing both methods leads to the agent not examining the cookbook, preventing it from receiving further instructions, which we can see from Figure 7(c) results in considerable performance loss, regressing to the baseline GATA.

## G.5 Code

Statement of intent. Upon publication of this work, we will release all code publicly on GitHub.

### G.5.1 Fixing the GATA code

We found two primary issues in the GATA code. First, we noticed that their implementation of the double Q-learning error was wrong. For Double Q-Learning, after performing some action  $a_t$  in state  $s_t$  and observing the immediate reward  $r_t$  and resulting state  $s_{t+1}$ , the Q-Learning error is defined per (Van Hasselt et al., 2016) as

$$Y_t = r_t + \gamma Q(s_{t+1}; \arg \max_a Q(s_{t+1}; a; \theta_t)) - Q(s_t; a_t; \theta_t) \quad (1)$$

where  $\theta_t$  and  $\theta_t^0$  are the parameters of the policy network and the target network, respectively. However, we noticed that the original code for GATA was computing the error as<sup>3</sup>

$$Y_t = r_t + r_{t+1} + \gamma Q(s_{t+1}; \arg \max_a Q(s_{t+1}; a; \theta_t)) - Q(s_t; a_t; \theta_t^0)$$

In other words, the reward for the stepped state was also being added to the error.

Second, we found that the double Q-learning error for terminal states was being incorrectly implemented. Specifically, when computing the error for the case where  $s_t$  is a terminal state, and therefore the stepped state  $s_{t+1}$  does not exist, the stepped state was not being masked. Additionally, presumably because of this initial error, terminal states were very rarely returned when sampling from experience, unless certain criteria were met.

We found fixing these issues improved GATA's performance considerably, which we demonstrated in Figure 4, and all our experimental results for GATA have this correction implemented.

## G.6 Training Curves

We present accompanying training curves for experiments reported in this work in Figure 8, Figure 9, Figure 10, Figure 11, and Figure 12. We report averaged curves of the normalized accumulated reward with bands representing the standard deviation.

## H Broader Impact

As (Adhikari et al., 2020) suggested, text-based games can be a proxy for studying human-machine

interaction through language. Human-machine interaction and relevant systems have many potential ethical, social, and safety concerns. Providing inaccurate policies or information or partially completing tasks in critical systems can have devastating consequences. For example, in health care, improper treatment can be fatal, or in travel planning, poor interactions can lose a client money.

Adhikari et al. (2020, section 7) identified several research objectives relating to language-based agents: improve the ability to make better decisions, allow for constraining decisions for safety purposes, and improve interpretability. We highlight how RL agents equipped with LTL instructions can improve in these areas. For constraining decisions, it may be desirable to do so in way that depends on the history, which LTL gives a way to keep track of. With respect to interpretability, we propose that monitoring the progression of instructions provides a mechanism for understanding where and when an agent might be making incorrect decisions, and provides the opportunity to revise instructions or attempt to fix the problem by other means.

Furthermore, we would like to suggest that works towards building better language agents should also emphasize the importance of completing instructions. To illustrate, for an agent to help a person half-way across a street, or to start but not finish a medical operation, may be worse than for it to do nothing at all. To that end, we have proposed using (game) success rate as a metric for future work, and demonstrated how LTL-GATA is very successful in the games it plays, relative to the state-of-the-art.

However, instruction following, especially overly literal instruction following, may not always be beneficial and can even be harmful. (Ammanabrolu et al., 2022) describe a good example where an agent in the Zork1 game breaks into a home and steals the items it needs. In that specific case, breaking into the home has no adverse effect on the agent's reward, and so it has no incentive not to perform this act. Violation of social norms like this are not modelled in our work, and can have negative impacts, even in less extreme cases. Furthermore, there are potential dangers of incorrect, immoral, or even misinterpreted instructions that lead to dangerous outcomes. Although we do not directly address these concerns in this work, they pose interesting directions for future work.

<sup>3</sup>[https://github.com/xingdi-eric-yuan/GATA-public/blob/c1afc3c9ab38256f839b3e0ddf8243796df5bd77/dqn\\_memory\\_prioritized\\_replay\\_buffer.py#L120-L123](https://github.com/xingdi-eric-yuan/GATA-public/blob/c1afc3c9ab38256f839b3e0ddf8243796df5bd77/dqn_memory_prioritized_replay_buffer.py#L120-L123)

<sup>4</sup><https://github.com/xingdi-eric-yuan/GATA-public/blob/c1afc3c9ab38256f839b3e0ddf8243796df5bd77/agent.py#L1353-L1369>

<sup>5</sup>[https://github.com/xingdi-eric-yuan/GATA-public/blob/c1afc3c9ab38256f839b3e0ddf8243796df5bd77/dqn\\_memory\\_prioritized\\_replay\\_buffer.py#L93-L102](https://github.com/xingdi-eric-yuan/GATA-public/blob/c1afc3c9ab38256f839b3e0ddf8243796df5bd77/dqn_memory_prioritized_replay_buffer.py#L93-L102)

Figure 8: Training curves of the normalized accumulated reward across various levels and on both the 20 (top) and 100 (bottom) game training sets. Bands represent the standard deviation. Note that on level 0, training curves for TDQN, GATA<sup>C</sup>, and GATA<sup>D</sup> were early stopped for achieving 0.95 normalized accumulated reward on the validation set for 5 episodes in a row.

Figure 9: Training curves (of normalized accumulated reward) for the study on LTL predicate format with single-token (Single) and multi-token (Multi) predicates. Bands represent the standard deviation.

Figure 10: Training curves (of normalized accumulated reward) for the comparison of GATA when trained with instructions (GATA<sup>D</sup>) versus when instructions are stripped from environment observations (GATA<sup>D</sup>-S). Agents were trained with 20 or 100 games, at increasing levels of task difficulty (level 1 vs level 2). Bands represent the standard deviation.

